

THE MAXIMUM ENTROPY ON THE MEAN METHOD, CORRELATIONS AND IMPLEMENTATION ISSUES

C. HEINRICH, J.-F. BERCHER AND G. DEMOMENT

Laboratoire des signaux et systèmes

Plateau de Moulon

91192 Gif-sur-Yvette Cedex, France [†]

Abstract. In this paper, we address the problem of solving linear inverse problems using the Maximum Entropy on the Mean Method (MEMM). The observation equation $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{n}$ links the observations \mathbf{z} and the supposed known degradation matrix \mathbf{H} to the object \mathbf{x} and the noise \mathbf{n} which have to be estimated. Given a reference measure ν defined jointly on \mathbf{x} and \mathbf{n} , the MEMM consists in selecting the distribution which is closest to ν according to the Kullback distance and whose mean satisfies the observation equation. The MEMM estimator $\hat{\mathbf{x}}$ is the mean of the selected distribution. This amounts to selecting $\hat{\mathbf{x}}$ as the minimizer of a convex cost function defined on \mathbf{x} and \mathbf{n} [1].

Up to now, this method yielded only separable objective functions. We present here an extension of the method, allowing penalization of linear functions of \mathbf{x} (such as differences between pixels for example) in the MEMM formalism. We also discuss algorithmic issues: the method amounts to minimizing a convex criterion and thus admits a dual formulation. We will take advantage of both formulations to reduce computational effort, which hadn't been achieved to this day since the problem had only been solved in the dual variables. We report two simulation examples, a deconvolution problem in spectroscopy and a Fourier synthesis example.

Key words: maximum entropy on the mean, maximum entropy, convex cost functions, duality, Legendre–Fenchel transform.

1. Introduction

We consider here linear inverse problems, whose observation model satisfies $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{n}$. We want to estimate the object \mathbf{x} and the noise \mathbf{n} , knowing the observations \mathbf{z} and the degradation matrix \mathbf{H} . Matrix \mathbf{H} will typically be a convolution or a Fourier matrix and vector \mathbf{x} a signal or an image. A typical resolution framework

[†]Email: heinrich@lss.supelec.fr

for this type of ill-posed problems is the minimization of a composite criterion

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}_{\mathbf{x}}(\mathbf{x}) + \mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x}),$$

yielding $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \mathcal{F}(\mathbf{x})$. This objective function comprises a regularization term $\mathcal{F}_{\mathbf{x}}$ and a term $\mathcal{F}_{\mathbf{n}}$ penalizing the residual noise $\mathbf{z} - \mathbf{H}\mathbf{x}$. Typical choices [2] are $\mathcal{F}_{\mathbf{x}}(\mathbf{x}) = \lambda_1 \sum_i \left| \frac{x_i}{\Delta_1} \right|^p + \lambda_2 \sum_j \left| \frac{x_{i+1} - x_i}{\Delta_2} \right|^q$ and $\mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x}) = \|\mathbf{z} - \mathbf{H}\mathbf{x}\|^2$.

Convexity of the cost function is a very important issue. Nonconvex cost functions are cumbersome to minimize and the estimator may not be a continuous function of the observations, continuity being a desirable property. Nevertheless, nonconvex cost functions may give better results than convex cost functions, especially regarding edge restoration. On the other hand, convex cost functions can be minimized using classical descent algorithms, such as gradient or conjugate gradient. Those functions guarantee continuity of the estimator. They are minimized at a reasonable computational cost.

Remembering Bayes' rule $p(\mathbf{x}|\mathbf{z}) \propto p(\mathbf{z}|\mathbf{x})p(\mathbf{x})$ linking the *a posteriori* distribution $p(\mathbf{x}|\mathbf{z})$ to the likelihood $p(\mathbf{z}|\mathbf{x})$ and the *a priori* distribution $p(\mathbf{x})$, $\hat{\mathbf{x}}$ can always be interpreted as a maximum *a posteriori* (MAP) estimator, providing the relations

$$\begin{aligned} p(\mathbf{x}) &\propto \exp(-\mathcal{F}_{\mathbf{x}}(\mathbf{x})), \\ p(\mathbf{z}|\mathbf{x}) &\propto \exp(-\mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x})) \end{aligned}$$

are satisfied.

We deal here with an alternative framework, the maximum entropy on the mean method, which allows to both construct and interpret convex cost functions. This class of convex cost functions comprises a lot of well known criteria, such as Shannon and Burg entropies and L_2 norm for example [3,1]. The method is closely linked to statistical mechanics. It was introduced by Navaza [4,5] to solve a inverse problem in crystallography. It has been further investigated by Dacunha-Castelle, Gamboa and Gassiat [3,6] from a mathematical point of view. Application of the method to inverse problems has been studied by Bercher, Le Besnerais and Demoment [7,1,8].

Given a reference measure ν defined on the objects to be estimated, this method consists in selecting the distribution \hat{p} minimizing the Kullback distance to ν among the distributions whose mean satisfies a given constraint. The MEM estimator is the mean of the selected distribution. The constraint is here an observation equation which will be assumed linear. Dealing with nonlinear problems is possible in the MEM formalism, but the optimization problem may no longer be convex.

We will assume that the constraints are qualified for all inverse problems considered here, which means that there exists a least one solution to any problem.

2. The maximum entropy on the mean method

To give a first sketch of the method, we will consider the problem of selecting an estimator $\hat{\mathbf{x}}$ among two candidates $\{\mathbf{x}_1, \mathbf{x}_2\}$. We will show that the MEM procedure amounts to selecting a distribution among an exponential family, and that deriving

the MEM estimator from this distribution is equivalent to minimizing a convex cost function defined on the object to be estimated.

Let

$$\mathcal{K}(p, \nu) = \begin{cases} \int \log \frac{dp}{d\nu} dp & \text{if } p \ll \nu, \\ +\infty & \text{otherwise,} \end{cases}$$

be the Kullback distance between the distribution p and the reference measure ν . Following the MEM principle, we derive:

$$\begin{aligned} \hat{p} &= \arg \min_{p \in \mathcal{C}} \mathcal{K}(p, \nu), \\ \mathcal{C} &= \mathcal{C}_{\mathbf{x}_1} \cup \mathcal{C}_{\mathbf{x}_2}, \\ \mathcal{C}_{\mathbf{x}_i} &= \{p \in \mathcal{M}_1^+ \mid E_p[\mathbf{u}] = \mathbf{x}_i\}, \\ \hat{\mathbf{x}} &= E_{\hat{p}}[\mathbf{u}], \end{aligned}$$

where \mathcal{M}_1^+ is the set of the probability measures.

Let $p_i = \arg \min_{p \in \mathcal{C}_{\mathbf{x}_i}} \mathcal{K}(p, \nu)$. We have $\hat{p} = \arg \min_{p \in \{p_1, p_2\}} \mathcal{K}(p, \nu)$, since

$$\forall p \in \mathcal{C} \quad \mathcal{K}(p, \nu) \geq \min(\mathcal{K}(p_1, \nu), \mathcal{K}(p_2, \nu)).$$

Selecting $\hat{\mathbf{x}}$ in $\{\mathbf{x}_1, \mathbf{x}_2\}$ thus amounts to comparing $\mathcal{F}(\mathbf{x}_1) \triangleq \mathcal{K}(p_1, \nu)$ to $\mathcal{F}(\mathbf{x}_2) \triangleq \mathcal{K}(p_2, \nu)$.

This can be extended to the choice of an estimator among any family \mathcal{E} (whether convex or not) of candidates.

We still have $\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{E}} \mathcal{F}(\mathbf{x})$, where $\mathcal{F}(\mathbf{x}) = \min_{\{p \mid E_p[\mathbf{u}] \in \mathcal{E}\}} \mathcal{K}(p, \nu)$. Convexity of $\mathcal{F}(\cdot)$ derives straightforwardly from convexity of $\mathcal{K}(\cdot, \nu)$.

In the large deviations theory, $\mathcal{F}(\cdot)$ is said to be a level-1 entropy, whereas $\mathcal{K}(\cdot, \nu)$ is the level-2 entropy (see [9] for developments on the large deviations theory and [10] for the links of the MEMM with statistical mechanics).

Let us now focus on distribution p_1 defined previously, so as to show that p_1 as well as \hat{p} belong to an exponential family. According to its definition, distribution p_1 minimizes the Kullback distance among all distributions of mean \mathbf{x}_1 . This convex constrained minimization problem is solved using the Lagrange multipliers method, which requires to find the stationary points of the Lagrangian

$$L = \mathcal{K}(p, \nu) + \boldsymbol{\mu}^t (\mathbf{x}_1 - E_p[\mathbf{u}]) + \mu_0 (1 - E_p[1])$$

associated to the problem. Vector $\boldsymbol{\mu}$ and scalar μ_0 are the Lagrange multipliers associated to the constraints. The Lagrangian will have to be minimized in variable p . We have:

$$L = \int \frac{dp}{d\nu} \left(\log \frac{dp}{d\nu} - \boldsymbol{\mu}^t \mathbf{u} - \mu_0 \right) d\nu + \boldsymbol{\mu}^t \mathbf{x}_1 + \mu_0.$$

Moreover, from

$$\forall y \in \mathbb{R}^{+*}, \forall \alpha \in \mathbb{R} \quad y(\log y - \alpha) \geq -\exp(\alpha - 1),$$

we derive

$$\frac{dp}{d\nu} \left(\log \frac{dp}{d\nu} - \boldsymbol{\mu}^t \mathbf{u} - \mu_0 \right) \geq -\exp(\boldsymbol{\mu}^t \mathbf{u} + \mu_0 - 1).$$

Hence

$$\forall p \in \mathcal{M}^+ \quad L \geq \mu_0 + \boldsymbol{\mu}^t \mathbf{x}_1 - \int \exp(\boldsymbol{\mu}^t \mathbf{u} + \mu_0 - 1) \, d\nu,$$

with equality if and only if

$$\frac{dp(\mathbf{u})}{d\nu(\mathbf{u})} = \exp(\boldsymbol{\mu}^t \mathbf{u} + \mu_0 - 1).$$

Let us define

$$\mathcal{F}^*(\boldsymbol{\mu}) \triangleq 1 - \mu_0 = \log \int \exp(\boldsymbol{\mu}^t \mathbf{u}) \, d\nu(\mathbf{u}),$$

assuring normalization of

$$dp(\mathbf{u}) = \exp(\boldsymbol{\mu}^t \mathbf{u} + \mu_0 - 1) \, d\nu(\mathbf{u}) = \exp(\boldsymbol{\mu}^t \mathbf{u} - \mathcal{F}^*(\boldsymbol{\mu})) \, d\nu(\mathbf{u}). \quad (1)$$

If we combine expressions of the Lagrangian L and of the exponential family (1), the Lagrangian reduces to $D(\boldsymbol{\mu}) \triangleq \boldsymbol{\mu}^t \mathbf{x}_1 - \mathcal{F}^*(\boldsymbol{\mu})$. Providing this concave function can be maximized, we define $\hat{\boldsymbol{\mu}}_1 = \arg \max_{\boldsymbol{\mu}} D(\boldsymbol{\mu})$. This defines the distribution \hat{p}_1 which is a stationary point of the Lagrangian L . By the way, we also established that $\mathcal{K}(\hat{p}_1, \nu) = \mathcal{F}(\mathbf{x}_1) = D(\hat{\boldsymbol{\mu}}_1)$ (dual attainment relation).

We derive at last equation $\mathcal{F}(\mathbf{x}) = \max_{\boldsymbol{\mu}} (\boldsymbol{\mu}^t \mathbf{x} - \mathcal{F}^*(\boldsymbol{\mu}))$, linking \mathcal{F} and \mathcal{F}^* as convex-conjugates or Legendre–Fenchel transforms of each other. Developments on the Legendre–Fenchel transform may be found in [9]. \mathcal{F}^* is said to be the log-partition function of ν , and \mathcal{F} is said to be the Cramér transform of ν . This uniquely defines \mathcal{F} if one is given a reference measure ν .

It must also be emphasized that the method wasn't exposed this way up to day. The former presentation (see [10] for example) restricted the application framework of the method, since it supposed that the observation equation was necessarily linear. This new presentation allows to deal with nonlinear problems and also allows to select an object among a discrete set of candidates. Of course, those are merely extensions enlightened by the presentation. These potentialities already existed before, but might have been hidden by the presentation.

3. The MEMM applied to linear inverse problems

In this section, we apply the MEMM reviewed in the previous section to the resolution of linear inverse problems. We will detail two equivalent formulations of the MEMM. We review additional developments needed to circumvent analytical difficulties otherwise seriously limiting the method.

We use notations previously established. The unknowns to be estimated are the object \mathbf{x} and the noise \mathbf{n} . The reference measure will have to be defined on \mathbf{x} and \mathbf{n} , and chosen such that $(\hat{\mathbf{x}}, \hat{\mathbf{n}})$ may be considered as the mean of a distribution near ν according to the Kullback distance. In particular, the selection procedure of \hat{p} implies that \hat{p} has the same support than ν . Since $\hat{\mathbf{x}}$ is the mean of \hat{p} , $\hat{\mathbf{x}}$ belongs to the convex hull of the support of ν . We take advantage of this property to constrain $\hat{\mathbf{x}}$ to belong to a given set. A reasonable choice is a separable reference measure $d\nu(\mathbf{u}, \mathbf{v}) = d\nu_{\mathbf{x}}(\mathbf{u}) \, d\nu_{\mathbf{n}}(\mathbf{v})$.

3.1. PRIMAL FORMULATION

We define here $d\nu(\mathbf{u}, \mathbf{v}) \triangleq d\nu_{\mathbf{x}}(\mathbf{u}) d\nu_{\mathbf{n}}(\mathbf{v})$ and $\boldsymbol{\mu} \triangleq [\boldsymbol{\mu}_{\mathbf{x}}^t, \boldsymbol{\mu}_{\mathbf{n}}^t]^t$, $\boldsymbol{\mu}_{\mathbf{x}}$ and \mathbf{u} having the same size as \mathbf{x} , $\boldsymbol{\mu}_{\mathbf{n}}$ and \mathbf{v} having the same size as \mathbf{n} . Following the reasoning exposed in the previous section, the log-partition function \mathcal{F}^* of the reference measure ν is given by:

$$\mathcal{F}^*(\boldsymbol{\mu}) = \log \int \exp(\boldsymbol{\mu}_{\mathbf{x}}^t \mathbf{u} + \boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{v}) d\nu(\mathbf{u}, \mathbf{v}).$$

We derive:

$$\begin{aligned} \mathcal{F}^*(\boldsymbol{\mu}) &= \log \int \exp(\boldsymbol{\mu}_{\mathbf{x}}^t \mathbf{u}) d\nu_{\mathbf{x}}(\mathbf{u}) + \log \int \exp(\boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{v}) d\nu_{\mathbf{n}}(\mathbf{v}), \\ &\triangleq \mathcal{F}_{\mathbf{x}}^*(\boldsymbol{\mu}_{\mathbf{x}}) + \mathcal{F}_{\mathbf{n}}^*(\boldsymbol{\mu}_{\mathbf{n}}). \end{aligned}$$

Hence we have:

$$\begin{aligned} \mathcal{F}(\mathbf{x}, \mathbf{n}) &= \max_{\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\mu}_{\mathbf{n}}} (\boldsymbol{\mu}_{\mathbf{x}}^t \mathbf{x} + \boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{n} - \mathcal{F}_{\mathbf{x}}^*(\boldsymbol{\mu}_{\mathbf{x}}) - \mathcal{F}_{\mathbf{n}}^*(\boldsymbol{\mu}_{\mathbf{n}})), \\ &= \max_{\boldsymbol{\mu}_{\mathbf{x}}} (\boldsymbol{\mu}_{\mathbf{x}}^t \mathbf{x} - \mathcal{F}_{\mathbf{x}}^*(\boldsymbol{\mu}_{\mathbf{x}})) + \max_{\boldsymbol{\mu}_{\mathbf{n}}} (\boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{n} - \mathcal{F}_{\mathbf{n}}^*(\boldsymbol{\mu}_{\mathbf{n}})), \\ &\triangleq \mathcal{F}_{\mathbf{x}}(\mathbf{x}) + \mathcal{F}_{\mathbf{n}}(\mathbf{n}). \end{aligned}$$

This shows that the cost function \mathcal{F} defined on \mathbf{x} and \mathbf{n} is separable if the reference measure is separable. The MEM estimator is selected by minimizing \mathcal{F} over all objects satisfying the observation equation, *i.e.* over all objects satisfying $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{n}$. This yields $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \mathcal{F}_{\mathbf{x}}(\mathbf{x}) + \mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x})$, which has the classical form of a regularized cost function. This is said to be the primal formulation of the MEM, and \mathcal{F} is called the primal cost function. We will now derive the dual formulation.

3.2. DUAL FORMULATION

To the minimization of a convex function (primal formulation) can be associated the maximization of a concave function (dual formulation). Both formulations yield the same solution. This relies on the Legendre–Fenchel duality theory [9]. We have:

$$\mathcal{F}_{\mathbf{x}}(\hat{\mathbf{x}}) + \mathcal{F}_{\mathbf{n}}(\hat{\mathbf{n}}) = \max_{\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\mu}_{\mathbf{n}}} (\boldsymbol{\mu}_{\mathbf{x}}^t \hat{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{n}}^t \hat{\mathbf{n}} - \mathcal{F}_{\mathbf{x}}^*(\boldsymbol{\mu}_{\mathbf{x}}) - \mathcal{F}_{\mathbf{n}}^*(\boldsymbol{\mu}_{\mathbf{n}})).$$

From $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \mathcal{F}_{\mathbf{x}}(\mathbf{x}) + \mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x})$, we derive:

$$\left. \frac{\partial \mathcal{F}_{\mathbf{x}}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\hat{\mathbf{x}}} - \mathbf{H}^t \left. \frac{\partial \mathcal{F}_{\mathbf{n}}(\mathbf{n})}{\partial \mathbf{n}} \right|_{\mathbf{n}=\mathbf{z}-\mathbf{H}\hat{\mathbf{x}}} = 0.$$

Hence $\hat{\boldsymbol{\mu}}_{\mathbf{x}} - \mathbf{H}^t \hat{\boldsymbol{\mu}}_{\mathbf{n}} = 0$ since, as a result of Fenchel duality,

$$\hat{\boldsymbol{\mu}}_{\mathbf{x}} = \left. \frac{\partial \mathcal{F}_{\mathbf{x}}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\hat{\mathbf{x}}} \quad \text{and} \quad \hat{\boldsymbol{\mu}}_{\mathbf{n}} = \left. \frac{\partial \mathcal{F}_{\mathbf{n}}(\mathbf{n})}{\partial \mathbf{n}} \right|_{\mathbf{n}=\hat{\mathbf{n}}}.$$

We have at last:

$$\mathcal{F}_x(\hat{\mathbf{x}}) + \mathcal{F}_n(\hat{\mathbf{n}}) = \max_{\boldsymbol{\mu}_n} (\boldsymbol{\mu}_n^t \mathbf{z} - \mathcal{F}_x^*(\mathbf{H}^t \boldsymbol{\mu}_n) - \mathcal{F}_n^*(\boldsymbol{\mu}_n)).$$

The primal formulation of the problem is thus linked to a dual formulation, which consists in maximizing the dual criterion $D(\boldsymbol{\mu}_n) \triangleq \boldsymbol{\mu}_n^t \mathbf{z} - \mathcal{F}_x^*(\mathbf{H}^t \boldsymbol{\mu}_n) - \mathcal{F}_n^*(\boldsymbol{\mu}_n)$. This amounts to selecting the parameters $(\hat{\boldsymbol{\mu}}_x = \mathbf{H}^t \hat{\boldsymbol{\mu}}_n, \hat{\boldsymbol{\mu}}_n)$ of the maximum entropy distribution belonging to the exponential family

$$dp(\mathbf{u}, \mathbf{v}) = \exp(\boldsymbol{\mu}_x^t \mathbf{u} + \boldsymbol{\mu}_n^t \mathbf{v} - \mathcal{F}_x^*(\boldsymbol{\mu}_x) - \mathcal{F}_n^*(\boldsymbol{\mu}_n)) d\nu_x(\mathbf{u}) d\nu_n(\mathbf{v}).$$

From $\hat{\boldsymbol{\mu}}_n$, we derive $\hat{\mathbf{x}}$ and $\hat{\mathbf{n}}$ using the so-called primal-dual relations:

$$\hat{\mathbf{x}} = \left. \frac{\partial \mathcal{F}_x^*(\boldsymbol{\mu}_x)}{\partial \boldsymbol{\mu}_x} \right|_{\boldsymbol{\mu}_x = \hat{\boldsymbol{\mu}}_x = \mathbf{H}^t \hat{\boldsymbol{\mu}}_n} \quad \text{and} \quad \hat{\mathbf{n}} = \left. \frac{\partial \mathcal{F}_n^*(\boldsymbol{\mu}_n)}{\partial \boldsymbol{\mu}_n} \right|_{\boldsymbol{\mu}_n = \hat{\boldsymbol{\mu}}_n}.$$

Remember also that $(\hat{\mathbf{x}}, \hat{\mathbf{n}})$ is the mean of the maximum entropy distribution. By the way, we can notice that $\forall \mathbf{x}, \mathbf{n}, \boldsymbol{\mu}_n \quad \mathcal{F}(\mathbf{x}, \mathbf{n}) \geq D(\boldsymbol{\mu}_n)$ and also that $\mathcal{F}(\hat{\mathbf{x}}, \hat{\mathbf{n}}) = D(\hat{\boldsymbol{\mu}}_n)$ (dual attainment relation).

3.3. FURTHER DEVELOPMENTS

Both formulations of the method require analytical expression of the log-partition function \mathcal{F}_x^* . It can't generally be computed if ν_x isn't separable (well known exceptions to this rule are non separable Gaussian and exponential reference measures).

In the sequel, reasoning will thus be restricted to separable reference measures. In a previously exposed computation, we proved that separable reference measures yield separable primal cost functions. We introduce an additional vector \mathbf{w} of reference measure ν_w so as to penalize interactions between pixels. If we choose all measures to be separable, this yields a complete separable primal cost function

$$\mathcal{F}(\mathbf{x}, \mathbf{w}, \mathbf{n}) = \mathcal{F}_x(\mathbf{x}) + \mathcal{F}_w(\mathbf{w}) + \mathcal{F}_n(\mathbf{n}).$$

This cost function will be minimized subject to $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{n}$ and $\mathbf{w} = \mathbf{Q}\mathbf{x}$, \mathbf{Q} being typically a differentiation matrix. We thus have:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \sum_i \mathcal{F}_{x_i}(x_i) + \sum_j \mathcal{F}_{w_j}((\mathbf{Q}\mathbf{x})_j) + \sum_k \mathcal{F}_{n_k}((\mathbf{z} - \mathbf{H}\mathbf{x})_k).$$

A straightforward computation similar to the one already detailed shows that the dual formulation of this problems writes:

$$\begin{cases} D(\boldsymbol{\mu}_w, \boldsymbol{\mu}_n) &= \boldsymbol{\mu}_n^t \mathbf{z} - \mathcal{F}_x^*(\mathbf{H}^t \boldsymbol{\mu}_n - \mathbf{Q}^t \boldsymbol{\mu}_w) - \mathcal{F}_w^*(\boldsymbol{\mu}_w) - \mathcal{F}_n^*(\boldsymbol{\mu}_n), \\ (\hat{\boldsymbol{\mu}}_w, \hat{\boldsymbol{\mu}}_n) &= \arg \max_{\boldsymbol{\mu}_w, \boldsymbol{\mu}_n} D(\boldsymbol{\mu}_w, \boldsymbol{\mu}_n), \\ \hat{\boldsymbol{\mu}}_x &= \mathbf{H}^t \hat{\boldsymbol{\mu}}_n - \mathbf{Q}^t \hat{\boldsymbol{\mu}}_w. \end{cases}$$

The MEM estimator is derived using

$$\hat{\mathbf{x}} = \left. \frac{\partial \mathcal{F}_{\mathbf{x}}^* (\boldsymbol{\mu}_{\mathbf{x}})}{\partial \boldsymbol{\mu}_{\mathbf{x}}} \right|_{\boldsymbol{\mu}_{\mathbf{x}} = \hat{\boldsymbol{\mu}}_{\mathbf{x}}}, \quad \hat{\mathbf{w}} = \left. \frac{\partial \mathcal{F}_{\mathbf{w}}^* (\boldsymbol{\mu}_{\mathbf{w}})}{\partial \boldsymbol{\mu}_{\mathbf{w}}} \right|_{\boldsymbol{\mu}_{\mathbf{w}} = \hat{\boldsymbol{\mu}}_{\mathbf{w}}}, \quad \hat{\mathbf{n}} = \left. \frac{\partial \mathcal{F}_{\mathbf{n}}^* (\boldsymbol{\mu}_{\mathbf{n}})}{\partial \boldsymbol{\mu}_{\mathbf{n}}} \right|_{\boldsymbol{\mu}_{\mathbf{n}} = \hat{\boldsymbol{\mu}}_{\mathbf{n}}}.$$

Selecting the dual variables $(\hat{\boldsymbol{\mu}}_{\mathbf{x}}, \hat{\boldsymbol{\mu}}_{\mathbf{w}}, \hat{\boldsymbol{\mu}}_{\mathbf{n}})$ amounts to selecting the maximum entropy distribution in the exponential family

$$dp(\mathbf{u}, \mathbf{s}, \mathbf{v}) = \exp(\boldsymbol{\mu}_{\mathbf{x}}^t \mathbf{u} + \boldsymbol{\mu}_{\mathbf{w}}^t \mathbf{s} + \boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{v} - \mathcal{F}_{\mathbf{x}}^* (\boldsymbol{\mu}_{\mathbf{x}}) - \mathcal{F}_{\mathbf{w}}^* (\boldsymbol{\mu}_{\mathbf{w}}) - \mathcal{F}_{\mathbf{n}}^* (\boldsymbol{\mu}_{\mathbf{n}})) \, d\nu_{\mathbf{x}}(\mathbf{u}) \, d\nu_{\mathbf{w}}(\mathbf{s}) \, d\nu_{\mathbf{n}}(\mathbf{v}).$$

We chose the reference measure $d\nu(\mathbf{u}, \mathbf{s}, \mathbf{v}) = d\nu_{\mathbf{x}}(\mathbf{u}) \, d\nu_{\mathbf{w}}(\mathbf{s}) \, d\nu_{\mathbf{n}}(\mathbf{v})$ to be separable. This isn't a contradiction since the constraint $\mathbf{w} = \mathbf{Q}\mathbf{x}$ applies to the mean of random variables and not to the realizations of the random variables. This constraint indeed links the parameters of the maximum entropy distribution by the relation $\hat{\boldsymbol{\mu}}_{\mathbf{x}} = \mathbf{H}^t \hat{\boldsymbol{\mu}}_{\mathbf{n}} - \mathbf{Q}^t \hat{\boldsymbol{\mu}}_{\mathbf{w}}$.

In a Bayesian framework, it appears that minimizing $\mathcal{F}_{\mathbf{x}}(\mathbf{x}) + \mathcal{F}_{\mathbf{w}}(\mathbf{Q}\mathbf{x}) + \mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x})$ is equivalent to a MAP estimation scheme with prior distribution $p(\mathbf{x}) \propto \exp(-\mathcal{F}_{\mathbf{x}}(\mathbf{x}) - \mathcal{F}_{\mathbf{w}}(\mathbf{Q}\mathbf{x}))$ and likelihood $p(\mathbf{z} | \mathbf{x}) \propto \exp(-\mathcal{F}_{\mathbf{n}}(\mathbf{z} - \mathbf{H}\mathbf{x}))$. The class of MEM estimators derived here is thus equivalent to the class of MAP estimators obtained with a given set of convex Markov random fields priors. Any MEM estimator is a MAP estimator and a MAP estimator is a MEM estimator providing there exists reference measures $\nu_{\mathbf{x}}, \nu_{\mathbf{w}}, \nu_{\mathbf{n}}$ verifying:

$$\begin{aligned} p(\mathbf{x}) &\propto \exp(-C(\nu_{\mathbf{x}}, \mathbf{x}) - C(\nu_{\mathbf{w}}, \mathbf{w})), \\ p(\mathbf{z} | \mathbf{x}) &\propto \exp(-C(\nu_{\mathbf{n}}, \mathbf{n})), \end{aligned}$$

where $C(p, \mathbf{y})$ is the Cramér transform of the distribution p , the Cramér transform being applied to \mathbf{y} .

4. Simulations

We report two simulation examples. We first deal with a deconvolution example in spectroscopy and then with a Fourier synthesis example.

Three hyperparameters k_a, k_b, k_c are introduced, the cost function to be minimized being $\mathcal{F}(\mathbf{x}, \mathbf{w}, \mathbf{n}) = k_a \mathcal{F}_{\mathbf{x}}(\mathbf{x}) + k_b \mathcal{F}_{\mathbf{w}}(\mathbf{w}) + k_c \mathcal{F}_{\mathbf{n}}(\mathbf{n})$, subject to $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{n}$ and $\mathbf{w} = \mathbf{Q}\mathbf{x}$. The corresponding dual function to be maximized is $\mathcal{F}^*(\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\mu}_{\mathbf{n}}) = k_c \boldsymbol{\mu}_{\mathbf{n}}^t \mathbf{z} - k_a \mathcal{F}_{\mathbf{x}}^*(\boldsymbol{\mu}_{\mathbf{x}}) - k_b \mathcal{F}_{\mathbf{w}}^*(\boldsymbol{\mu}_{\mathbf{w}}) - k_c \mathcal{F}_{\mathbf{n}}^*(\boldsymbol{\mu}_{\mathbf{n}})$ subject to $\boldsymbol{\mu}_{\mathbf{x}} = (k_c \mathbf{H}^t \boldsymbol{\mu}_{\mathbf{n}} - k_b \mathbf{Q}^t \boldsymbol{\mu}_{\mathbf{w}}) / k_a$. The hyperparameters are chosen empirically.

We tried several optimization schemes out (primal and dual objective functions optimized with pseudo-conjugate gradient descent method and pixel by pixel (Gauss-Seidel) descent method). It seems hard to derive a definite general rule from the simulations we carried out. In any case, it must be checked if optimization has been properly completed. In the primal domain, one will thus have to check if

$$\mathcal{F}(\hat{\mathbf{x}}, \hat{\mathbf{w}}, \hat{\mathbf{n}}) = \mathcal{F}^* \left(\left. \frac{\partial \mathcal{F}_{\mathbf{x}}(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x} = \hat{\mathbf{x}}}, \left. \frac{\partial \mathcal{F}_{\mathbf{w}}(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w} = \hat{\mathbf{w}}}, \left. \frac{\partial \mathcal{F}_{\mathbf{n}}(\mathbf{n})}{\partial \mathbf{n}} \right|_{\mathbf{n} = \hat{\mathbf{n}}} \right).$$

A similar relation has to be checked in the dual domain.

Moreover, one should be careful not to stop iterations before the optimum is reached when operating in the dual variables. This would yield unsatisfactory results, much more unsatisfactory than when stopping iterations before the optimum is reached in the primal variables.

Simulations we carried out speak in favor of pixel by pixel minimization. On the one hand, given a pixel to be updated, every function evaluation is achieved at a very low computational cost (instead of computing $\mathbf{H}\mathbf{x}$ or $\mathbf{H}^t\boldsymbol{\mu}_n$, we use updating formulas). Gradient or pseudo-conjugate gradient descent requires computation of $\mathbf{H}\mathbf{x}$ or $\mathbf{H}^t\boldsymbol{\mu}_n$ for any function evaluation. On the other hand, separation of the variables in the pixel by pixel scheme allows more efficient optimization of the function (figure 2c).

4.1. DECONVOLUTION

We used a mixture of two Gamma laws as a reference measure for the pixels and a Gaussian law as a reference measure for the noise.

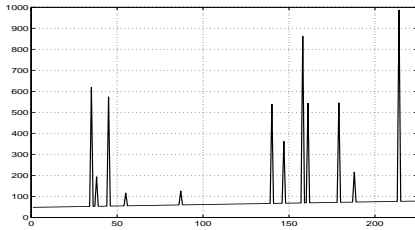


Fig. 1a: original signal

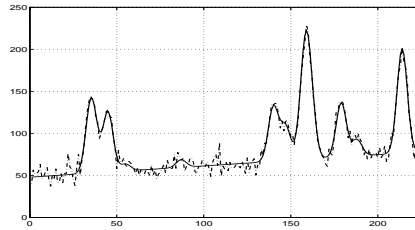


Fig. 1b: observations; SNR : 16 dB

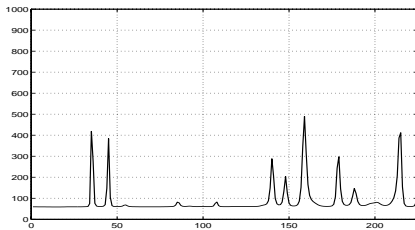


Fig. 1c: mem restauration

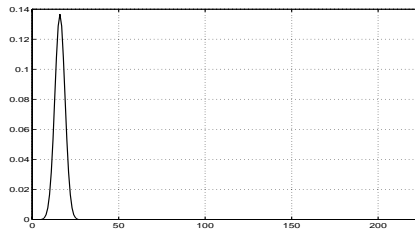


Fig. 1d: convolution kernel

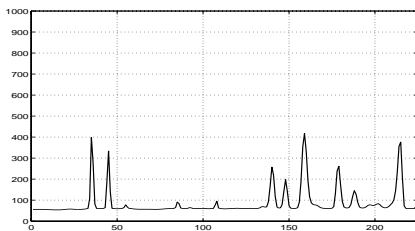


Fig. 1e: Lp restauration (p=1.1)

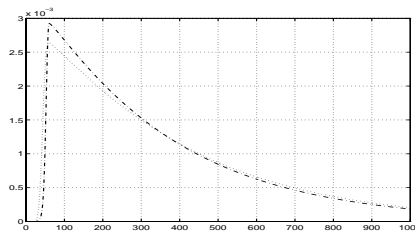


Fig. 1f: priors

Figure 1f (priors) depicts $\exp(-x^p)$ (dotted line) and $\exp(-\mathcal{F}(x))$ (dashdotted line). Those functions would play the role of a prior distribution in a Bayesian

framework.

4.2. FOURIER SYNTHESIS

This example is taken from [11]. The data are the real parts of the twenty first Fourier coefficients which are supposed to be perfectly known. We used a mixture of two exponential laws (defined on \mathbb{R}^+) as a reference measure for the pixels, a mixture of two symmetric exponential laws as a reference measure for the difference between neighboring pixels and a Gaussian law as a reference measure for the noise.

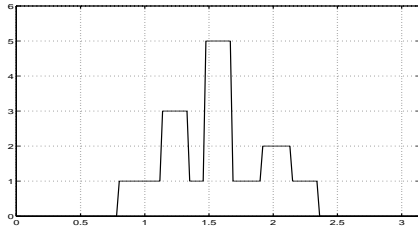


Fig. 2a: original signal

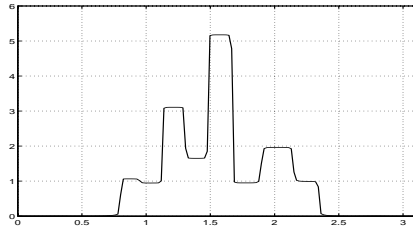


Fig. 2b: mem reconstruction

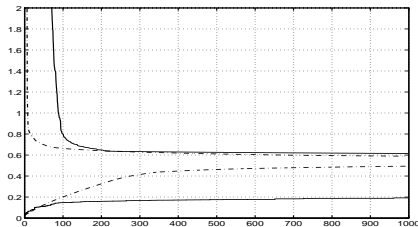


Fig. 2c: criteria vs cpu time (seconds)

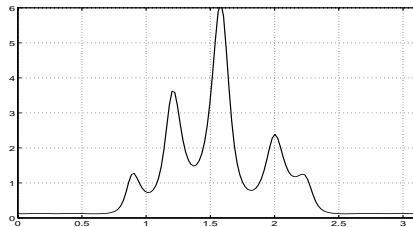


Fig. 2d: reconstruction; seperable criterion

Figure 2c depicts the value of the objective functions (above: primal domain; below: dual domain; solid lines: pseudo-conjugate gradient descent; dashdotted lines: pixel by pixel descent). Figure 2d depicts reconstruction obtained without penalizing differences between pixels, which would be comparable to the results reported in [11].

5. Conclusion

This communication deals with both theoretical and algorithmic issues. On the one hand, we account for penalization of interactions between pixels in the MEM formalism. This introduces additional terms in the criterion, which nevertheless remains strictly convex. On the other hand, we focus on implementation issues. This speaks in favor of pixel by pixel minimization (Gauss-Seidel strategy), either in the primal or in the dual domain. It appears in the case investigated here that minimization in the primal domain is less expensive than in the dual domain. This doesn't hold as a general and definite conclusion. It merely shows that many strategies are available and that the choice of the suitable one is certainly problem dependant.

Additional work on the method is under way. Other approaches for penalizing interactions in the MEMM have been proposed: Urban [12] decomposes the error n_i as $n_i = n_i^0 + \sum_j n_{ij}$, where $n_{ij} = n_{ji}$ and uses a separable reference measure on $\{n_i^0, n_{ij}\}$. Moreover, study of mathematical aspects of the method by Csiszár, Gamboa and Gassiat is in progress.

The MEM provides a framework for constructing convex criteria. This framework also allows to interpret well known criteria. But, in the opinion of the authors, this framework is by no mean exclusive. In particular, [2] suggests many interesting convex and nonconvex criteria in a Bayesian approach.

References

1. J.-F. Bercher, G. Le Besnerais, and G. Demoment, *The maximum entropy on the mean method, noise and sensitivity*. Maximum entropy and Bayesian methods, Kluwer Academic Publishers, Cambridge, U.K., 1995.
2. C. Bouman and K. Sauer, "A generalized Gaussian image model for edge-preserving MAP estimation," *IEEE Transactions on Image Processing*, **IP-2**, pp. 296–310, July 1993.
3. F. Gamboa, *Méthode du maximum d'entropie sur la moyenne et applications*. PhD thesis, Université de Paris-Sud, Orsay, December 1989.
4. J. Navaza, "On the maximum entropy estimate of electron density function," *Acta Crystallographica*, **A-41**, pp. 232–244, 1985.
5. J. Navaza, "The use of non-local constraints in maximum-entropy electron density reconstruction," *Acta Crystallographica*, **A-42**, pp. 212–223, 1986.
6. D. Dacunha-Castelle and F. Gamboa, "Maximum d'entropie et problème des moments," *Annales de l'Institut Henri Poincaré*, **26**, (4), pp. 567–596, 1990.
7. G. Le Besnerais, *Méthode du maximum d'entropie sur la moyenne, critères de reconstruction d'image, et synthèse d'ouverture en radio-astronomie*. PhD thesis, Université de Paris-Sud, Orsay, December 1993.
8. J.-F. Bercher, *Développement de critères de nature entropique pour la résolution des problèmes inverses linéaires*. PhD thesis, Université de Paris-Sud, Orsay, February 1995.
9. R. S. Ellis, *Entropy, Large Deviations, and Statistical Mechanics*, Springer-Verlag, New York, 1985.
10. C. Heinrich, J.-F. Bercher, G. Le Besnerais, and G. Demoment, "Restoration of spiky signals: a new optimal estimate and a comparison," in *Proceedings of IEEE ICASSP*, (Detroit, U.S.A.), pp. 877–880, May 1995.
11. L. K. Jones and V. Trutzer, "Computationally feasible high-resolution minimum-distance procedures which extend the maximum-entropy method," *Inverse Problems*, **5**, pp. 749–766, 1989.
12. B. Urban, "Retrieval of atmospheric thermodynamical parameters using satellite measurements with a maximum entropy method," tech. rep., Centre National de Recherches Météorologiques, Toulouse, France, April 1996.